

5

10

15

20

25

10

15

4. A document classification method for classifying a document based on contents of the document of which contents contains a plurality of items, said document classification method comprising the steps of:

designating at least one of the items
contained in the document input in the inputting step;

classifying the document by using the
10 converted data produced in the converting step.

25

6. The document classification system as claimed in claim 4, wherein the converting step includes the step of inserting a predetermined sign between sets of data corresponding to the items so as to facilitate separation of each data corresponding to each item in the converted data.

10

547
A17 7. A processor readable medium storing program code causing a computer to classify a document based on contents of the document of which contents contains a plurality of items, comprising:

15

first program code means for inputting document data corresponding to the document data;

second program code means for designating at least one of the items contained in the document;

20

third program code means for converting the document data into converted data so that the converted data contains only data corresponding to the item designated by the second program code means; and

25

fourth program code means for classifying the document by using the converted data produced by the third program code means.

8. The processor readable medium as claimed
in claim 7, wherein the fourth program code means
includes fifth program code means for producing a
feature vector representing a feature of the converted
5 data so as to classify the document in accordance with
the feature vector.

10

9. The processor readable medium as claimed
in claim 7, wherein the third program code means
includes sixth program code means for inserting a
predetermined sign between sets of data corresponding to
15 the items so as to facilitate separation of each data
corresponding to each item in the converted data.

20

10. A document classification system for
classifying a document according to contents of the
document, said document classification system
comprising:
25 input means for inputting document data of the

analyzing means for analyzing the document
data so as to obtain analysis information; 3

transforming function calculating means for calculating a representation transforming function used for projecting the document feature vector onto a space in which similarity between the document feature vectors is reflected;

15 classification means for classifying the document based on similarity between the document feature vectors transformed by the vector transforming means; and

25

11. The document classification system as
claimed in claim 10, further comprising inner product
calculating means for calculating an inner product
between the document feature vectors, wherein said
5 representation transforming function calculating means
calculates the representation transforming function by
using the inner product.

10

12. The document classification system as
claimed in claim 11, further comprising document
similarity information setting means for setting
15 document similarity setting information including data
representing an author of the document and a date of
production of the document, wherein said representation
transforming function calculating means calculates the
representation transforming function by using the inner
20 product and the document similarity information.

25

vector storing means for storing the document
feature vector produced by said vector producing means;
5 and

10

20

15. The document classification system as claimed in claim 14, further comprising transforming function correcting means for correcting the representation transforming function calculated by said transforming function calculating means when the feature dimension is changed due to a correction of the document feature vector by said vector correcting means so that the document feature vector is transformed by said vector transforming means in accordance with the changed feature dimension.

16. The document classification system as claimed in claim 10, further comprising:
transforming function correction instructing means for sending an instruction regarding a process to be applied on a feature dimension of the representation transforming function; and

transforming function correcting means for correcting the representation transforming function based on a content of the instruction sent from said transforming function correction instructing means.

25

5

15

25

10 21. The document classification system as
claimed in claim 10, further comprising:
an initial cluster centroid designating means
for designating an initial cluster centroid; and
initial cluster centroid registering means for
15 registering the initial cluster centroid designated by
said initial cluster centroid designating means,
wherein said classification means classifies
the document in accordance with the initial cluster
centroid registered by said initial cluster centroid
20 registering means.

22. The document classification system as
25 claimed in claim 21, wherein the initial cluster

5

15

20

25. The document classification system as
25 claimed in claim 21, wherein the initial cluster

5

26. A document classification method for classifying a document according to contents of the document, said document classification method comprising the steps of:

inputting document data of the document;

analyzing the document data so as to obtain analysis information;

producing a document feature vector with respect to the document data based on the analysis information;

calculating a representation transforming function used for projecting the document feature vector onto a space in which similarity between the document feature vectors is reflected;

transforming the document feature vector by using the representation transforming function;

classifying the document based on similarity between the document feature vectors transformed in the

5

10

15

25

5

10

15

25

10

25

5

15

25

5

10

15

20

25

39. The document classification method as claimed in claim 37, wherein the initial cluster centroid designated in the step of designating is the document feature vector.

5

40. The document classification method as claimed in claim 37, wherein the initial cluster centroid designated in the step of designating is the analysis information obtained in the step of analyzing.

15

41. The document classification method as claimed in claim 37, wherein the initial cluster centroid designated in the step of designating is the result of classification stored in the step of storing.

25

42. A processor readable medium storing program code causing a computer to classify a document according to contents of the document, comprising:

second program code means for analyzing the document data so as to obtain analysis information;

fourth program code means for calculating a representation transforming function used for projecting the document feature vector onto a space in which similarity between the document feature vectors is reflected;

fifth program code means for transforming the document feature vector by using the representation transforming function;

sixth program code means for classifying the
20 document based on similarity between the document
feature vectors transformed by the fifth program code
means; and

seventh program code means for storing a
result of classification performed by the classification
25 means.

10

15

20

25

document feature vector produced by the third program code means; and

eleventh program code means for storing the representation transforming function calculated by the
5 fourth program code means.

10 46. The processor readable medium as claimed in claim 42, further comprising twelfth program code means for correcting the document feature vector before the document feature vector is transformed by the fifth
15 program code means, a correction being performed by processing one of the document feature vector and a feature dimension constituting the document feature vector in accordance with a rule established by characteristics of words extracted by the second program
20 code means.

47. The processor readable medium as claimed in claim 46, further comprising thirteenth program code
25 means for correcting the representation transforming

10

15

20

25

Adh
A7

25